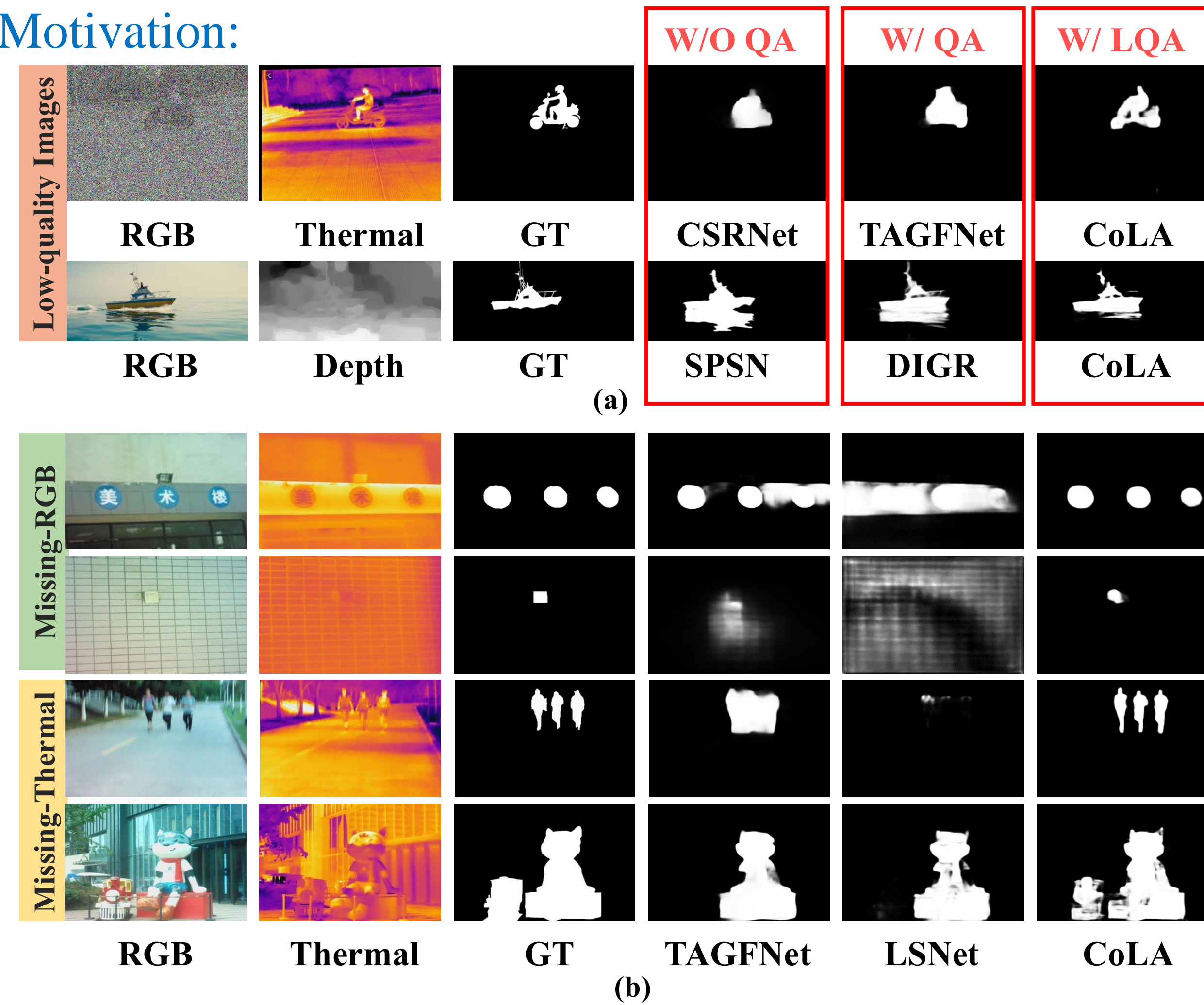
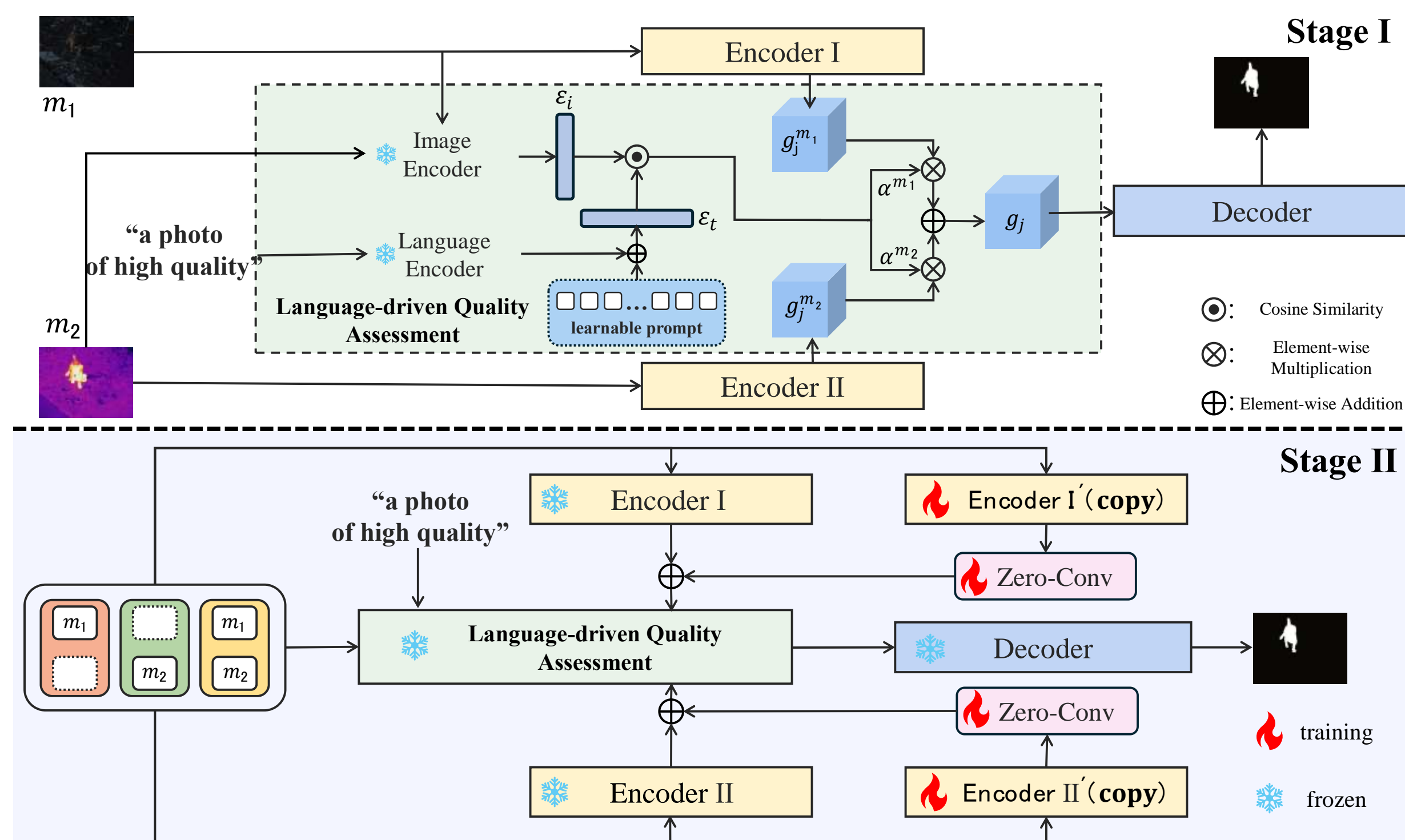


Motivation:



Our Method:

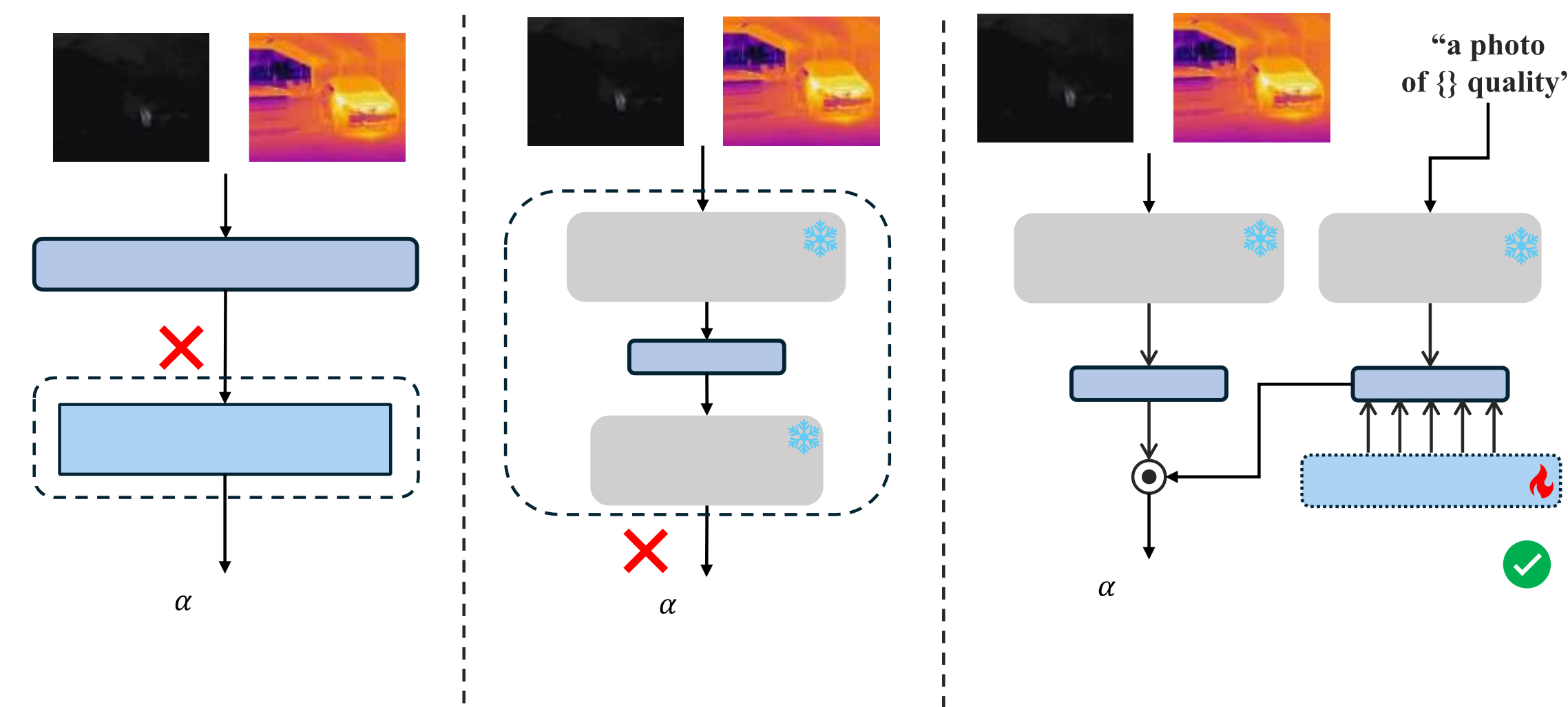


We propose the CoLA framework, which includes:

- 1) **Language-driven Quality Assessment (LQA)** to recalibrate image in put using a vision-language model, reducing noise impact without extra annotations,
- 2) **Conditional Dropout (CD)** to improve model performance when modalities are missing.

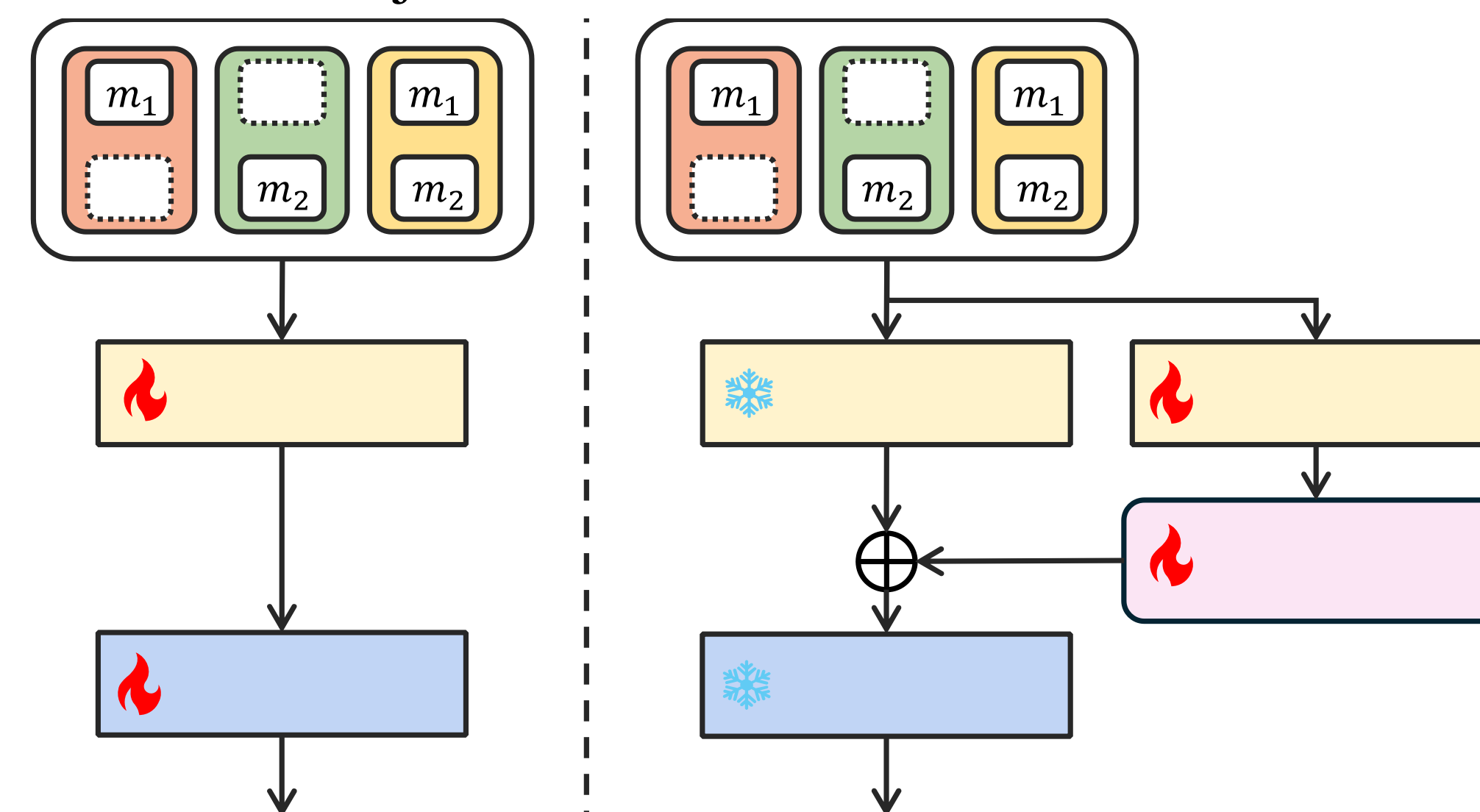
Language-driven Quality Assessment:

Architectural comparison of (a) No-Reference Method, (b) Pre-trained Assessment and (c) Our LQA.



Conditional Dropout:

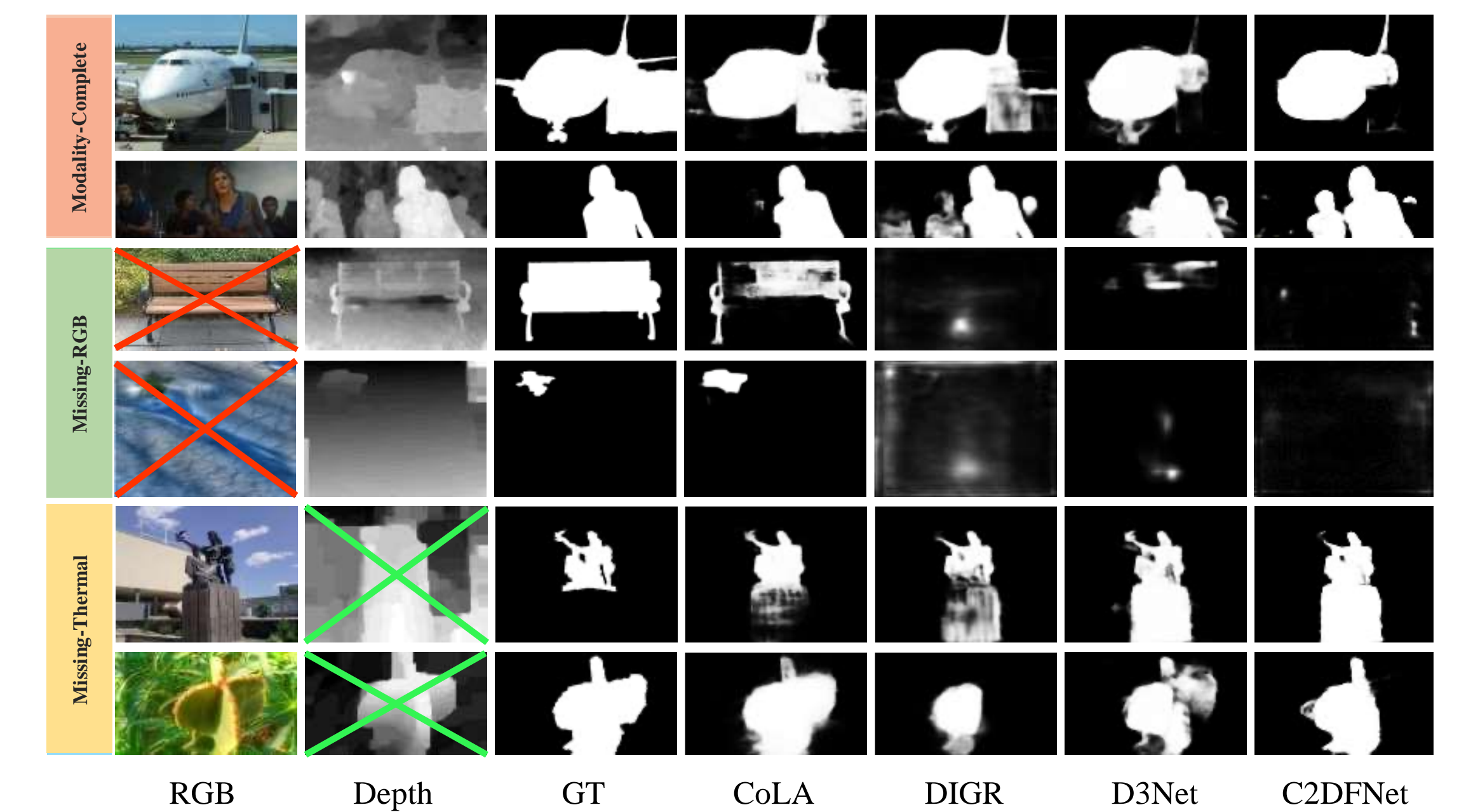
Architectural comparison of (a) modality dropout and (b) Conditional Dropout in dual-modal object detection.



Results:

| RGB-T SOD | | | | RGB-D SOD | | | |
|----------------|------------|-----------|---------------|-----------------------------|----------|---------|-----------|
| Methods | VT821 Full | VT821 RGB | VT821 Thermal | Methods | DES Full | DES RGB | DES Depth |
| DCNet (2022) | 0.841 | 0.644 | 0.78 | C ² DFNet (2022) | 0.940 | 0.568 | 0.902 |
| TAGFNet (2023) | 0.825 | 0.727 | 0.771 | SPSN (2022) | 0.950 | 0.793 | 0.908 |
| LSNet (2023) | 0.829 | 0.687 | 0.749 | HiDANet (2023) | 0.979 | 0.907 | 0.922 |
| Ours | 0.849 | 0.752 | 0.817 | Ours | 0.963 | 0.947 | 0.926 |

Visualization:



Code Available



Wechat



WhatsApp